

# Robust approximation of the conditional mean for applications of Machine Learning

Amy Parkes<sup>a,\*</sup>, Josef Camilleri<sup>b</sup>, Dominic Hudson<sup>a</sup>, Adam Sobey<sup>a,c</sup>

<sup>a</sup>*Maritime Engineering, University of Southampton, Southampton, SO17 1BJ, UK*

<sup>b</sup>*Silverstream Technologies Ltd, 1 St Vincent Street, London, W1U 4DA, UK*

<sup>c</sup>*Data-centric Engineering, The Alan Turing Institute, The British Library, NW1 2DB, London, UK*

---

## Abstract

Machine Learning approaches are increasingly used in a range of applications. They are shown to produce low conventional errors but in many real applications fail to model the underlying input-output relationships. This is because the error measures used only predict the conditional mean under some restrictive assumptions, often not met by the data we extract from applications. However, new approaches to Machine Learning, for example using Evolutionary Computation, allow a range of alternative error measures to be used. This paper explores the use of the Fit to Median Error measure in machine learning regression automation, using evolutionary computation in order to improve the approximation of the ground truth. When used alongside conventional error measures it improves the robustness of the learnt input-output relationships to the conditional median. It is compared to traditional regularisers to illustrate that the use of the Fit to Median Error produces regression neural networks which model more consistent input-output relationships. The problem considered is ship power prediction using a fuel-saving air lubrication system, which is highly stochastic in nature. The networks optimised for their Fit to Median Error are shown to approximate the ground truth more consistently, without sacrificing conventional Minkowski-r error values.

**Keywords:** Machine Learning, Genetic Algorithms, Neural Architecture Search, Optimisation, Application

---

---

\*Corresponding author  
Preprint submitted to *Artificial Intelligence* (Amy Parkes)

## 7 1. Robust approximation of the ground truth

8 Machine learning regression models are increasingly being used in indus-  
9 trial and engineering contexts for high-stakes decision making, automation and  
10 control. These methods often produce low conventional error values, yet only  
11 produce physically inconsistent results under a number of assumptions[3], which  
12 are often not met in real applications. Outside of these constraints they cannot  
13 generalise off test set and so cannot be relied upon to model the ground truth  
14 of the system. In real-world applications, where the output can have a direct ef-  
15 fect on human life or the environment, model accuracy alone is not sufficient. It  
16 has been demonstrated that for many applications minimising traditional error  
17 measures cannot guarantee an accurate approximation of the ground truth [28].  
18 This is due to a poor inductive bias, the inherent prioritisation of one solution  
19 over another [2], produced by conventional error measures which are based on  
20 Minkowski-r metrics [10].

21 This trust can be increased by manually tuning to remove overfitting or to  
22 provide a solution that makes more sense to the user. However, the expert  
23 knowledge and domain experience required to properly tune a machine learn-  
24 ing method manually are not always available in industry. Genetic algorithms  
25 are also increasingly used to search a method’s hyperparameter space more ef-  
26 ficiently [29] [15] [1]; which minimise conventional error measures on a test set,  
27 often combined with lowering the complexity of the network. This automation  
28 exacerbates the lack of interpretability, as models have a large flexibility, and  
29 prediction accuracy is prioritised. A low conventional error is often achieved  
30 without certainty that the method has modelled the correct internal functions.  
31 Regularisation hyperparameters can be optimised alongside other neural net-  
32 work parameters [26] [19], which increases the search space and creates more  
33 flexibility for methods to produce ‘accurate’ predictions and avoid overfitting.  
34 These approaches improve the modelling of the ground truth in scenarios ad-  
35 hering to the assumptions in the proof in [3], which are

- 36 1. the datapoints are independent;

2. the distribution of the target variable is to be deterministic of the input with Gaussian noise, e.g.  $y = \phi(x) + \epsilon$  where  $\phi$  depends only on input variable  $x$ , and  $\epsilon \sim N(0, \sigma^2)$ ;
3. the standard deviation of noise,  $\sigma$ , is not dependent on the input  $x$ ;
4. and the data set and neural network must be sufficiently large.

and under which minimum Minkowski-r error values approximate the conditional average of the dataset. This is because the inductive bias from the loss function guides the input-output relationships towards the conditional average, while the regularisation stops overfitting by simplifying the input-output relationships being modelled. However, these assumptions are restrictive and it is noted that few regression applications adhere to them. For example, one assumption is that the dataset is homoscedastic. In scenarios not adhering to these assumptions, network regularisation simplifies the relationships being modelled but this does not necessarily improve the generality, or model the ground truth.

The Fit to Median Error measure [22] produces models with a better fit to the true input-output relationships, when used in conjunction with conventional error measures. The use of this error measure removes the need for assumptions 2. and 3. to hold to model the ground truth. Which is achieved by regularising the learnt input-output relationships to the conditional median of the training dataset: the median output value, conditioned on each isolated input variable in turn [3]. For many regression applications the conditional medians are a good approximation of the ground truth input-output relationships, and therefore biasing the input-output relations learnt by the regression method towards the conditional medians produces models with more robust input-output relationships. The Fit to Median Error measure has been shown to produce models with a better approximation of the true input-output relations by trialling it on an artificial dataset, where the input-output relations can be fully defined, the results of this study can be found in [22]. As yet, the Fit to Median error it has not been explored as part of an automated approach.

This paper therefore explores the automation of neural network training to a

new problem, with a focus on producing a network which accurately models the ground truth. To achieve this, the cMLSGA multi-objective genetic algorithm is used to tune the hyperparameters of neural networks. The study compares the ground truth representation of a neural network when a genetic algorithm optimises the network’s hyperparameters to reduce the Mean Fit to Median Error measure and compares it to standard regularization using l1, l2 and dropout, and to a network optimised to minimise the Maximum Absolute Error. It is illustrated that neural network regularisation methods (l1, l2 and dropout) can be replaced by the use of the Mean Fit to Median performance measure as an objective in the genetic algorithm, reducing the complexity of the search space and producing networks which more consistently model the ground truth.

## 2. Neural Networks Parameters

A challenging regression problem is ship power prediction for a vessel using air lubrication to reduce fuel consumption. It is chosen to be used in this study as it violates the assumptions in [3], where the noise in the output space is non-Gaussian and heteroscedastic, it is also likely that there is not enough data. In this situation, correctly modelling the ground truth and accurate prediction is required but there is limited understanding of the ground truth [21], meaning physics-informed approaches are not applicable. The literature shows that shaft powering of a vessel can be predicted with average accuracies of between 1.5-5% error with the use of a regression neural network trained with high frequency data from the vessel [23], [24], [16], [12] and [18]. All neural network applications to ship power prediction in the literature use a combination of local searches and domain knowledge to identify hyperparameter values. The addition of an air lubrication device increases the complexity of the regression problem, as the system interacts with a number of interrelated input variables.

Previous applications of neural networks to ship power prediction use between 1 and 3 hidden layers [17] [20], and between 5 and 300 neurons in each hidden layer [12]. To provide a sufficiently large search space to allow verification, or otherwise, of these parameters a maximum of 4 hidden layers and 1000

neurons in each layer are used. The majority of the literature treats the problem as time-invariant and use feed-forward networks, so no recurrent parameters are optimised. As the optimiser or activation functions are rarely documented in the literature, the state-of-the-art optimisers and activation functions available in the Keras framework [4] are used in the optimisation, Table 1.

Table 1: Selected Neural Network Hyperparameters

| Hyperparameter          | Value or set  |
|-------------------------|---|
| Layers                  | [1,4]   |
| Neurons in each layer   | [1,1000]  |
| Epochs                  | Increasing from 1-20 for increasing generations   |
| Early stopping patience | 5   |
| Loss function           | Mean Absolute Error   |
| Performance measures    | Mean Absolute Relative Error, Maximum Absolute Relative Error, Mean Fit to Median Error |
| Optimiser               | SGD, Adam [14], Nadam [5], RMSprop [11], Adagrad [6], Adadelata [30], Adamax [14]       |
| Activation function     | ReLU, sigmoid, softmax, softplus, softsign, tanh, selu, elu                             |
| l1 & l2 Rates           | 0, 0.01,0.001,0.0001,0.00001  |
| Dropout                 | [0,0.9)   |
| Initialiser             | Random Normal ( $\mu = 0, \sigma = 0.1$ )   |

The number of epochs and early stopping procedure are not optimised, as there was a need for predictable compute requirements and allowing the optimisation of these parameters leads to unpredictable run times. The number of epochs to train each network increases for increasing generation number in the genetic algorithm, from 1 epoch in the first 15 generations to 20 in the final 15. This was also implemented to reduce compute and it was validated that when more than 20 epochs were allowed, that the early stopping, with a patience of 5, stopped the training within 20 epochs for the majority of networks. The loss function is similarly not optimised, the Mean Absolute Error is used, as the conditional medians are closer to the ground truth input-output relationships in these datasets than the conditional means.

113 The performance measures, or the genetic algorithm’s fitness functions, are  
114 the Mean Absolute Relative Error, the Maximum Absolute Relative Error and  
115 the Mean Fit to Median Error. Different combinations of these, alongside the  
116 use of L1, L2<sup>1</sup> and dropout<sup>2</sup> regularisation parameters in the search space are  
117 compared to illustrate the effect of different types of regularisation.

### 118 3. cMLSGA Parameters

Table 2: Selected cMLSGA Hyperparameters

| Hyperparameter                   | Value or set        |
|----------------------------------|---------------------|
| Algorithm at Individual Level    | HEIA, IBEA          |
| Crossover Type & Rate            | SBX & DE, 1         |
| Mutation Type & Rate             | Polynomial,<br>0.08 |
| Number of eliminated collectives | 1                   |
| Generations between elimination  | 10                  |
| Population size                  | 1000                |
| Generations                      | 300                 |
| Proportion elite                 | 10%                 |

119 In this study cMLSGA<sup>3</sup> is selected as it shows the top performance on a  
120 range of evolutionary benchmarking problems [9] and practical problems [8].  
121 Genetic algorithms are increasing used to tune neural network hyperparame-  
122 ters including regularisation parameters for use on new problems [13]. Many  
123 approaches have multiple genetic algorithm objectives, although these all min-  
124 imise an error measure and a measure of network complexity [27] and [25]. The  
125 use of multiple different performance measures as objectives is yet to be explored  
126 in the literature.

127 Four approaches are investigated in this study, summarised in Table 3, for  
128 approach (GAi) and (GAii) the genetic algorithm cMLSGA optimises all vari-

<sup>1</sup>L1 and L2 encourage simple relationships to be modelled by penalising large weights. L1 adds the absolute values of the weights to the cost function and L2 adds the squared values of the weights to the cost function.

<sup>2</sup>Dropout works by turning off a certain proportion of neurons randomly at each run of the network, and it’s use is shown to be equivalent to using a Bayesian network [7]

<sup>3</sup>The code for cMLSGA is available at <https://github.com/12yuens2/cmlsga-jmetalpy>.

Table 3: Genetic Algorithm Approaches

| Approach     | Objective(s)                                    | Network Regularisation |
|--------------|---|------------------------|
| <b>GAi</b>   | Mean Absolute Error                             | l1, l2 and dropout     |
| <b>GAii</b>  | Mean Absolute Error<br>Maximum Absolute Error   | l1, l2 and dropout     |
| <b>GAiii</b> | Mean Fit to Median Error<br>Mean Absolute Error | None                   |
| <b>GAiv</b>  | Mean Absolute Error<br>Maximum Absolute Error   | None                   |

ables in Table 2, including the l1 and l2 regularisation rate and the dropout rate of the networks. Although it is advised that l2 regularisation and dropout are not used in the same network the genetic algorithms are provided with zero options for all regularisation parameters, to identify if one is preferable in this scenario.

Approach (**GAi**) is a single objective genetic algorithm optimising the Mean Absolute Error which is compared to a multi-objective formulation where the (**GAii**) approach optimises both Mean Absolute Error and Maximum Absolute Error. For approaches (**GAiii**) and (**GAiv**) no network regularisation parameters are optimised: l1, l2 and dropout rates are all set permanently to zero. They avoid producing networks that have overfitted by the use of two performance metrics as multi-objectives, (**GAiii**) uses the Mean Fit to Median and Mean Absolute Errors to be minimised and (**GAiv**) uses the Maximum Absolute and Mean Absolute. All approaches use 40 CPUs with 2.0 GHz Intel Skylake processors and 192 GB of DDR4 memory, and take less than 3 days, this setup may not be feasible for widespread industrial application, although it is suggested it is within reach of some industries.

#### 4. Data

The data used in this study are from a large vessel equipped with the Silverstream® Air Lubrication System. The air lubrication system works through use

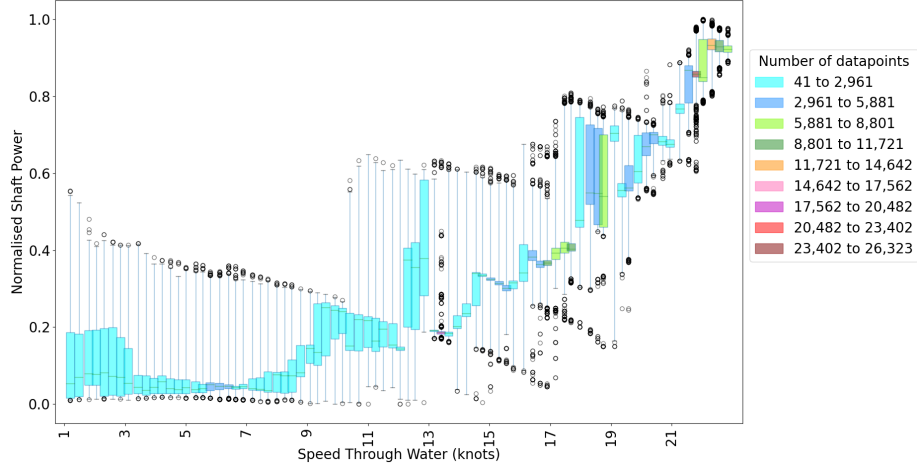


Figure 1: The distribution of the observed shaft powers for half knot bins of speed through the water for dataset where the system is off. In the box and whisker plots the boxes contain 50% of the distribution and the whiskers extend to the datum which is at 1.5 times the interquartile range.

149 of fluid sheering to create an air microbubble carpet directly captured within  
 150 the boundary layer on the ship hull bottom. The bubble carpet reduces the  
 151 frictional resistance thereby increasing the speed and reducing the shaft power.  
 152 Compressors provide a constant supply of air to the hull bottom to maintain  
 153 a uniform bubble carpet operated at the optimal compressor power that max-  
 154 imises the energy balance. The study is performed on both system on and  
 155 system off datasets, however for brevity only results for system off are presented  
 156 as they show similar performance. This prediction is required for a baseline  
 157 determination of how the system is working, but the relationships between the  
 158 power, weather, ocean and operating conditions are complex and difficult to  
 159 model.

160 The variables considered in this study are the shaft power, speed through  
 161 water, relative wind speed and direction, draught and trim, with shaft power  
 162 the target variable. These are selected based on a detailed study into variable  
 163 selection for shaft power prediction [20]. The speed through water is selected  
 164 over the speed over ground, for use as an input variable, as it is more hydrody-



manically relevant and its accuracy is validated by comparison to the speed over ground. The dataset is cleaned by removing rows with missing or non-physical values and all datapoints below 0.05 normalised shaft power are removed. The dataset is split into two using the air lubrication system status: system on and system off, where system on is defined as air lubrication system power greater than zero. The system on dataset contains 352,690 datapoints and system off contains 237,962. The data is split into training, testing and validation sets of 70%, 15% and 15% respectively. Each network in the genetic algorithm trains on a randomly sampled 35,000 datapoints from the training set and uses randomly sampled sets of size 7,500 from validation and testing sets for validation during training, and testing to produce the fitness of the network for the genetic algorithm. The errors stated in the paper are from networks on the Pareto fronts of each approach, which are validated on the full testing set.

The datasets contain large regions of sparse data in all input variable domains, this is exemplified by the ship speed domain where each half-knot interval below 16 knots contains less than 0.8% of the data, which accounts for more than half the speed domain, Figure 1. In addition, the boxplot ranges and outliers show high heteroscedicity with idiosyncratic noise caused by situations where the angle of the propeller blades is varied to achieve the required speed. This highlights the complexity in developing models of the powering of this vessel, as the dataset also contains the effects from other latent variables, such as piloting behaviour and route taken.

## 5. Optimisation including regularisation parameters: (GAi) and (GAii)

Previous studies predicting ship powering using neural networks report that l1, l2 and elastic net increase both test set and off-test set errors and that optimal values for both l1 and l2 are zero. Therefore the genetic algorithm setup is biased towards low and zero values of regularisation rates by using a set of exponentially decreasing values and an explicit zero option.

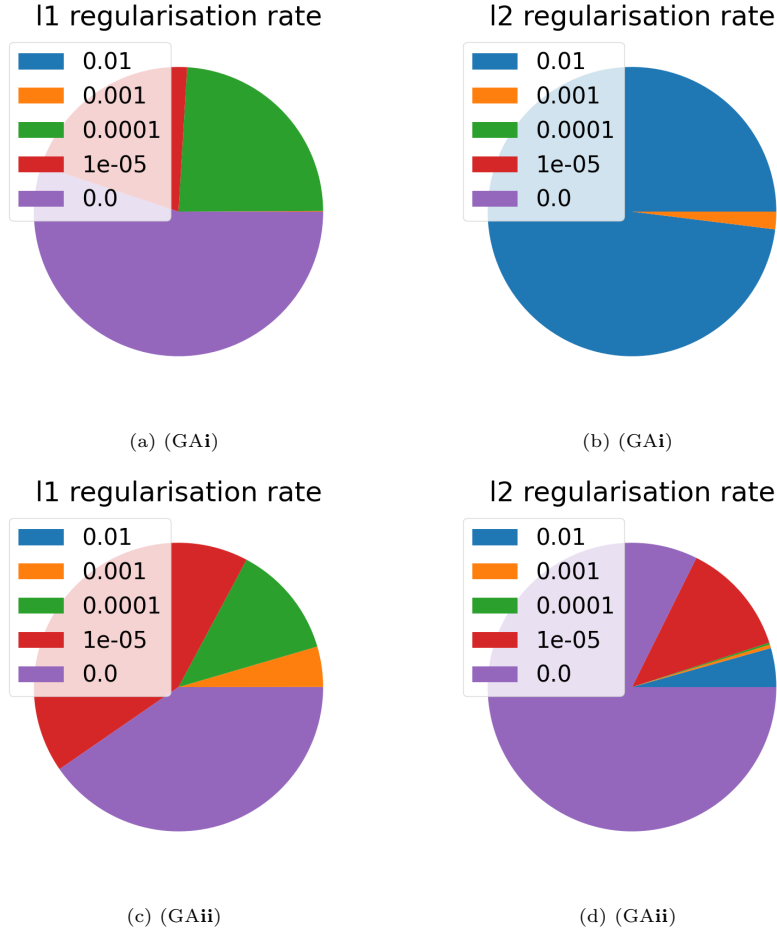
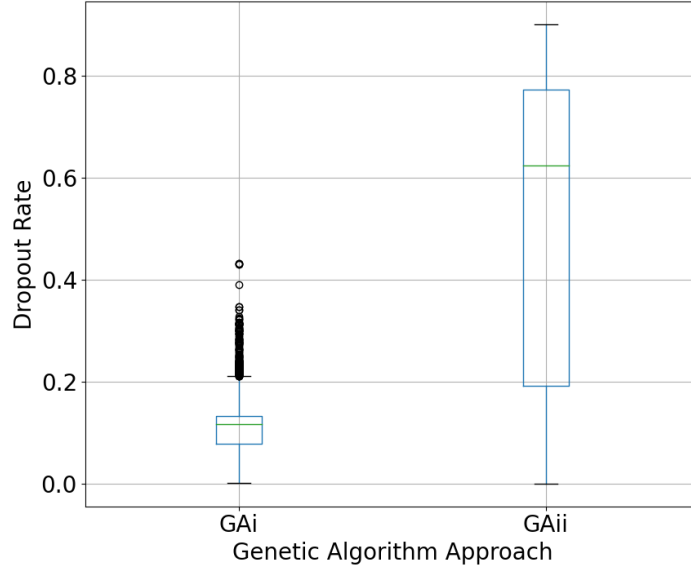
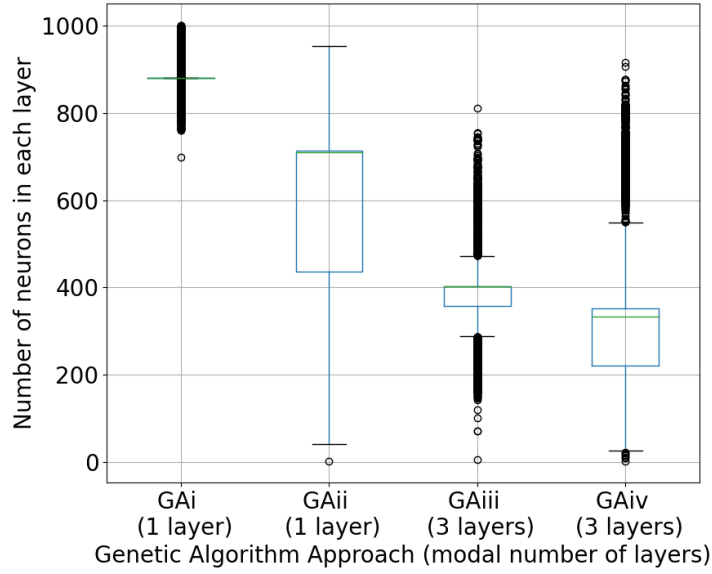


Figure 2: Distribution of regularisation rates for networks in the last 15 generations of (GAi) cMLSGA with multi-objectives of minimising Maximum and Mean Absolute Error for (a) l1 and (b) l2 and (GAii) cMLSGA with the single objective of minimising Mean Absolute Error for (c) l1 and (d) l2

193 The single objective (GAi) fails to identify that zero regularisation rates  
 194 produce the lowest errors, favouring networks with the highest possible rate of  
 195 l2 (0.01), Figure 2b. (GAi) produces networks with the highest Mean Absolute  
 196 Relative Errors of all the approaches,  $(5.19 \pm 0.00)\%$  from Figure 4a. In contrast,  
 197 (GAii) favours lower l1 and l2 rates of 0 or 0.00001, Figures 2c and 2d, which  
 198 results in networks with the lowest Mean Absolute Relative Errors of all four  
 199 approaches, on average, with a value of  $(2.87 \pm 0.45)\%$ , shown in Figure 4a.



(a) a



(b) b

Figure 3: (a) Dropout rate for networks in the last 15 generations of cMLSGA with (GAi) the single objective of minimising Mean Absolute Error and (GAii) multi-objectives of minimising Maximum and Mean Absolute Error and (b) the number of neurons in each layer for networks in the last 15 generations of cMLSGA with (GAi), (GAii), (GAiii) and (GAiv).

200 This is around 0.5% higher than the lowest documented error for ship power  
201 prediction.

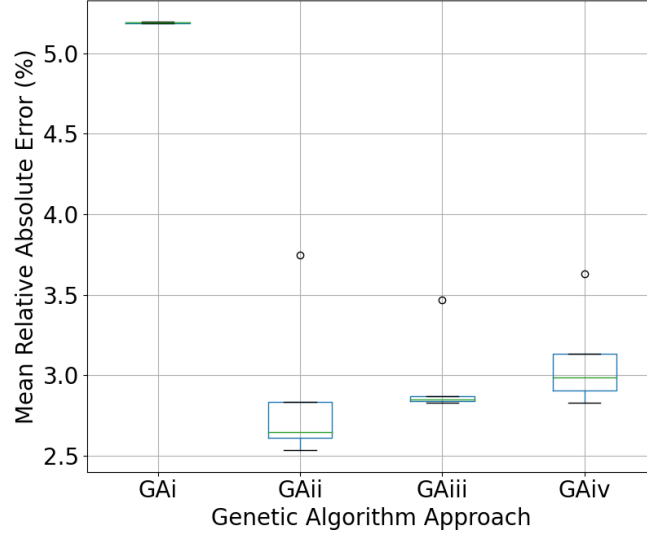
202 It is posited that the high error for the single objective problem is directly  
203 related to the use of a large l2 regularisation rate, as noted in previous studies  
204 for ship power prediction. It is possible that the use of a multi-objective search  
205 algorithm for a single-objective problem means that the optimal hyperparameters  
206 can't be found, resulting in large errors. The implementation also requires  
207 restrictions in the number of epochs used for training in the initial generations,  
208 it is possible this biases (GAi) towards certain size networks, where higher l2  
209 rates are preferable. This hypothesis is supported by the fact that 74.4% of  
210 networks in the first 15 generations of (GAi) have 1 hidden layer, and that  
211 over 99.8% of the networks in the final 15 generations have 1 hidden layer, with  
212  $880 \pm 17$  neurons in this layer, Figure 3b. This is significantly more neurons than  
213 those in the hidden layer of networks in the final 15 generations of (GAii) which  
214 range from 3-952 with a median value of 709, Figure 3b. The added objective of  
215 minimising Maximum Absolute Error in (GAii) may cause these slightly smaller  
216 networks to be more attractive as they are in a sense regularised by their size,  
217 as they have reduced modelling flexibility therefore are less likely to overfit and  
218 produce high Maximum Absolute Errors.

219 Another explanation for the difference in l2 rates chosen by (GAi) and (GAii)  
220 is the equivalence of l2 and dropout. Since l2 and dropout are equivalent up  
221 to a Fisher transformation, their use in conjunction is not recommended. The  
222 evidence for this is that (GAi) favours the highest l2 rate and has a median  
223 dropout rate in the final 15 generations of 0.116, whereas (GAii) favours the  
224 zero l2 rate and has a median dropout rate of 0.624, Figure 3a. This illustrates  
225 that the genetic algorithms will chose either l2 or dropout to minimise the Mean  
226 Absolute Relative Error. The l1 rates also support this hypothesis, as chosen  
227 rates for l1 regularisation in the final 15 generations are more comparable for  
228 (GAi) and (GAii).

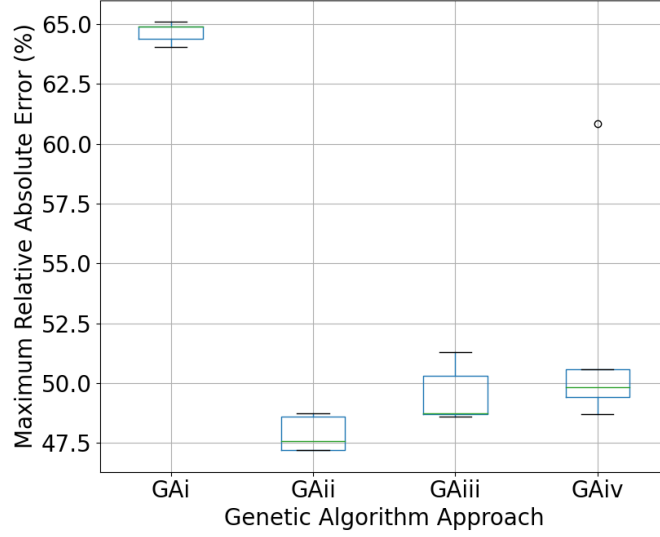
## 229 6. Optimisation using multiple performance measures: (GAiii) and 230 (GAiv)

231 For approaches (GAiii) and (GAiv) all neural network regularisation param-  
232 eters are set to zero. The regularisation is performed by minimising different  
233 network performance measures, the Mean Absolute and Mean Fit to Median for  
234 (GAiii), and the Mean Absolute and Maximum Absolute for (GAiv). The trade-  
235 off between the two objectives produces regularised neural networks, without  
236 explicitly changing the architecture or loss function. The Mean Fit to Median is  
237 chosen as it indicates how close the relationships modelled by a network are to  
238 the conditional averages of the dataset, in many regression examples this is akin  
239 to the ground truth input-output relationships [22]. The Maximum Absolute  
240 is chosen as for many industrial applications of machine learning the maximum  
241 prediction error is more pertinent than the mean error. The Mean Absolute  
242 Error is used instead of the Mean Squared Error in both approaches, as the  
243 conditional medians are closer to the ground truth input-output relationships  
244 in these datasets than the conditional means.

245 Differently shaped networks are favoured by (GAiii) and (GAiv), compared  
246 to (GAi) and (GAii), focusing on networks with 3 hidden layers and on average  
247 less than 400 neurons in each layer, Figure 3b. These networks have 51 times  
248 the number of connections than the networks chosen in (GAi) and (GAii).  
249 Apart from (GAi), (GAiii) has the most consistently sized networks in the  
250 final 15 generations, with an interquartile range of 46 neurons, compared to  
251 (GAiv) which have an interquartile range of 131 neurons. It is suggested that  
252 as the Mean Fit to Median Error biases networks towards specific input-output  
253 relationships, there is a smaller range of potential network architectures which  
254 habitually model these relationships. Whereas networks which minimise the  
255 Maximum Absolute Error are less restricted and can model a wider range of  
256 input and output relationships.



(a) a



(b) b

Figure 4: (a) Mean Relative Absolute Error and (b) Maximum Absolute Error from cMLSGA with (GAi) the single objective of minimising Mean Absolute Error and (GAii) multi-objectives of minimising Maximum and Mean Absolute Error, both optimising the parameters for l1, l2 regularisation and dropout in the networks, and (GAiii) and (GAiv) which do not use network regularisation but minimise Mean Fit to Median and Maximum Absolute Error respectively, alongside Mean Absolute Error

257 The Mean Absolute Relative Errors from networks in the Pareto fronts  
 258 are  $(2.97 \pm 0.25)\%$  for (GAiii) and  $(3.10 \pm 0.28)\%$  for (GAiv). It is expected  
 259 that (GAiv) would produce higher Mean Absolute Relative Errors as discussed  
 260 above, minimising the Maximum Absolute Error should bias predictions to-  
 261 wards the midpoint of the conditional output distributions, whereas minimising  
 262 the Mean Absolute Error should bias predictions towards the median of these  
 263 distributions. As it is established that noise in the output distribution is non-  
 264 Gaussian, Figure 1, these values will not align so some sacrifice in Mean Absolute  
 265 Error is expected from (GAiv). Both (GAiii) and (GAiv) produce comparable  
 266 Maximum Absolute Errors, of  $(49.5 \pm 1.1)\%$  and  $(51.9 \pm 4.5)\%$ . It is suggested  
 267 that this is because, although the conditional median output value and condi-  
 268 tional midpoint output value do not align for the majority of the input domain,  
 269 they are sufficiently close to produce comparable Maximum Absolute Errors.

270 Across all four approaches, the genetic algorithm producing networks with  
 271 the highest Mean Absolute Error is the approach which does not provide extra  
 272 weighting to sparse areas of data. The approaches minimising Maximum Abso-  
 273 lute Error are implicitly biased away from networks which predict the majority  
 274 of the testing datapoints correctly, but predict one datapoint poorly, favouring  
 275 networks which predict all testing datapoints to a moderate degree of error.  
 276 Approach (GAiii) more explicitly weights prediction in sparse areas of data by  
 277 favouring networks which model the conditional median of the dataset across all  
 278 input domains, irrespective of the quantity of data across each input domain.  
 279 The regression problem of ship power prediction is chosen in part because of it's  
 280 irregular data distribution; more than 9% of the dataset lies in less than a 0.5  
 281 knot interval of ship speed, Figure 1. This provides an explanation for the high  
 282 testing errors from (GAi), where only the Mean Absolute Error is minimised,  
 283 there is little incentive for the genetic algorithm to produce networks which  
 284 generalise across the full range of the input domain well.

## 285 7. Comparison of the ground truth approximation

286 Five networks selected from the four different approaches the learnt rela-  
 287 tionship between an input, the ship speed, and the output, shaft power, for the  
 288 networks in the Pareto front of each approach are visualised, Figure 5. These  
 289 are extracted with the following procedure: set all but one input variable to  
 290 be constant at the mode; cycle the remaining variable from its minimum to its  
 291 maximum recorded values with 150 points evenly spaced along the domain and  
 292 run the new dataset through the trained network.

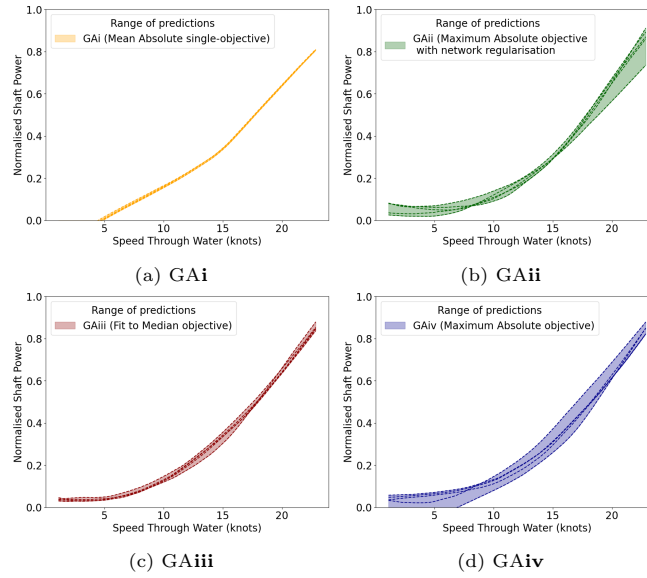


Figure 5: The learnt speed-power curves from 5 networks on the Pareto fronts of (GAii), (GAiii), (GAiv) and the 5 networks producing lowest Mean Absolute Relative Error from (GAi).

293 The approach which produces the most consistent speed-power relationships  
 294 is (GAi), with an average variation of 1.8%<sup>4</sup>, Figure 5a. However, the relation-  
 295 ship modelled by the 5 networks with the lowest Mean Absolute Relative Error

---

<sup>4</sup>Average variation in just the speed-power curves are discussed in this section, but it is verified that all input-power curves follow the same trends with variation around 0.5% across input variables.



296 in (GAi) all approximate a piece-wise linear relationship which clearly underfits  
 297 the dataset in Figure 1. The expected trend between ship speed through the  
 298 water and shaft power is a cubic polynomial, therefore as well as producing the  
 299 highest Mean Absolute Relative Errors, networks chosen by (GAi) model the  
 300 ground truth input-output relationships the worst out of the four approaches.  
 301 Both (GAii) and (GAiv) produce 5 fairly consistent speed-power curves, with  
 302 average variations of 5.9% and 10% respectively, Figures 5b and 5d. Both ap-  
 303 proaches approximate smooth polynomial curves, although the degrees of the  
 304 polynomials might differ, as multiple curves intersect at various points along  
 305 the speed axis. The spread of learnt relationships is greater at the highest and  
 306 lowest speeds for (GAii), with a decrease in spread for speeds of around 15  
 307 knots, where many of the curves intersect. The curves from (GAiv) show equal  
 308 spread across the speed domain.

309 The approach with both accurate and consistent learnt speed-power curves  
 310 is (GAiii), with limited intersections of curves and an average spread of 3.0%.  
 311 It is suggested that the reason using the Mean Fit to Median Error as an ob-  
 312 jective in a multi-objective genetic algorithm produces more consistent learnt  
 313 relationships, is because instead of encouraging the networks to model more  
 314 simple relationships it encourages the networks to model the conditional me-  
 315 dian functions of the dataset, supported by the increase in network connections.  
 316 Whereas the other approaches leave room for networks to fail to model the  
 317 conditional averages, especially in irregularly distributed and non-normally dis-  
 318 tributed datasets. The Mean Absolute Error values from networks selected by  
 319 (GAiii) are on average 0.1% higher than those from (GAii), and the Maximum  
 320 Absolute Error values are 1.6% higher.

321 A limitation of the approach is that the Fit to Median Error measure will  
 322 perform best at improving fit to the ground truth on datasets which violate  
 323 the assumptions stated in section 1; the ship powering example is chosen to  
 324 illustrate this as it provides a clearly heteroscedastic dataset with a cubic speed-  
 325 power relationship. For applications where noise profiles are Gaussian, and  
 326 there are no latent or interrelated input variables, the Fit to Median Error will

not improve the fit to ground truth but will perform the same as conventional Minkowski-r metrics, either Mean Squared or Mean Absolute Error depending on the convexity of input-output relationships.

Interestingly, the approach producing the lowest Mean Absolute and Maximum Absolute Errors does not model the ground truth the most accurately. This creates a potential for negative societal impacts, as the standard performance metrics for regression neural networks do not provide a full picture of performance or expected behaviour. Accurate input-output relationships are essential for safe applications of machine learning in the real world, especially when automated methods are used to replace experienced professionals. (GAii) demonstrates the same accuracy of approach as those with standard network regularisation, but with a better fit to the ground truth. This approach bypasses the need to use, and therefore to optimise the parameters of the regularisation methods. If evolutionary computation is already being used to optimise network parameters, then compute is saved by removing the network regularisation parameters l1, l2 and dropout. (GAiii) completed 300 generations in 46hours whereas (GAii) required 12 hours more computation to complete 300 generations.

## 8. Conclusion

The error measures used in our current Machine Learning approaches can provide an accurate point-wise estimate, without being able to approximate the input-output relationships in many real world datasets. New error measures are applied that match to the expected mean under less restrictive conditions. Three different approaches are compared: one to minimise the Maximum Absolute Error of the networks, which includes standard regularization using l1, l2 and dropout; and two which do not use any network regularisation, one minimising the Mean Fit to Median and one to minimise the Maximum Absolute Error. The results show that all three approaches give similar Mean Absolute Errors from networks on their Pareto fronts, from 2.9% for the approach with regularisation to 3.1% for the approach minimising Maximum Absolute Error.

357 However, using a new measure, the Mean Fit to the Median, a considerably  
358 better approximation of the expected mean can be made, with a spread in pre-  
359 dicted input-output curves of 3% compared to a spread of 6% for the approach  
360 using regularisation and 10% when minimising the Maximum Absolute Error  
361 and where it is the only approach where the input-output curves don't cross.

## 362 Acknowledgments

363 The work was also kindly supported by the Lloyds Register Foundation.  
364 The authors acknowledge the use of the IRIDIS High Performance Computing  
365 Facility, and associated support services at the University of Southampton, in  
366 the completion of this work. The authors would also like to thank Silverstream  
367 Technologies Ltd for providing the dataset used in this study and their support.  
368 We would also like to thank the Knowledge Transfer Network for sponsoring  
369 this research under grant KTP012306.

## 370 References

- 371 [1] Alexandrov, I.A., K. A. K. V. C. L. [2023], *High Tech and Innovation*  
372 *Journal* **4**.
- 373 [2] Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zam-  
374 baldi, V., Malinowski, M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner,  
375 R., Gulcehre, C., Song, F., Ballard, A., Gilmer, J., Dahl, G., Vaswani, A.,  
376 Allen, K., Nash, C., Langston, V., Dyer, C., Heess, N., Wierstra, D., Kohli,  
377 P., Botvinick, M., Vinyals, O., Li, Y. and Pascanu, R. [2018], 'Relational  
378 inductive biases, deep learning, and graph networks'.  
379 **URL:** <https://arxiv.org/abs/1806.01261>
- 380 [3] Bishop, C. [1995], *Neural Networks for Pattern Recognition*, Oxford Uni-  
381 versity Press, chapter 6, pp. 194–225.
- 382 [4] Chollet, F. et al. [2015], 'Keras', <https://keras.io>.

- [5] Dozat, T. [2016], ‘Incorporating nesterov momentum into adam’, *International Conference on Learning Representations 2016* .
- [6] Duchi, J., Hazan, E. and Singer, Y. [2011], ‘Adaptive subgradient methods for online learning and stochastic optimization’, *Journal of Machine Learning Research* **12**(61), 2121–2159.
- [7] Gal, Y. and Ghahramani, Z. [2016], Dropout as a bayesian approximation: representing model uncertainty in deep learning, *in* ‘International Conference on Machine Learning, New York, USA’.
- [8] Grudniewski, P. A. and Sobey, A. J. [2019], Do general genetic algorithms provide benefits when solving real problems?, *in* ‘2019 IEEE Congress on Evolutionary Computation (CEC)’, IEEE, pp. 1822–1829.
- [9] Grudniewski, P. A. and Sobey, A. J. [2021], ‘cMLSGA: a co-evolutionary multi-level selection genetic algorithm for multi-objective optimization’.  
**URL:** <https://arxiv.org/abs/2104.11072>
- [10] Hanson, S. and Burr, D. [1987], Minkowski-r back-propagation: learning in connectionist models with non-euclidian error signals., *in* ‘Neural Information Processing Systems (NIPS 1987)’.
- [11] Hinton, G., Srivastava, N. and Swersky, K. [2012], ‘Lecture 6a:overview of mini-batch gradient descent’, [http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture\\_slides\\_lec6.pdf](http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf).
- [12] Jeon, M., Noh, Y., Shin, Y., Lim, O., Lee, I. and Cho, D. [2018], ‘Prediction of ship fuel consumption by using an artificial neural network’, *Journal of Mechanical Science and Technology* **32**(12), 5785–5796.
- [13] Jin, Y., Okabe, T. and Sendhoff, B. [2004], Neural network regularization and ensembling using multi-objective evolutionary algorithms, *in* ‘Proceedings of the 2004 Congress on Evolutionary Computation (IEEE Cat. No.04TH8753)’, Vol. 1.

- [14] Kingma, D. P. and Ba, J. [2014], ‘Adam: A method for stochastic optimization’.  
**URL:** <https://arxiv.org/abs/1412.6980>
- [15] Kumar, P., Batra, S. and Raman, B. [2021], ‘Deep neural network hyperparameter tuning through twofold genetic approach’, *Soft Computing* pp. 1–25.
- [16] Le, L., Lee, G., Park, K. and Kim, H. [2020], ‘Neural network-based fuel consumption estimation for container ships in korea’, *Maritime Policy & Management* pp. 1–18.
- [17] Leifsson, L., Sævarsdóttir, H., Sigursson, S. and Vésteinsson, A. [2008], ‘Grey-box modeling of an ocean vessel for operational optimization’, *Simulation Modelling Practice and Theory* **16**(8), 923–932.
- [18] Liang, Y., Niu, D. and Hong, W. [2019], ‘Short term load forecasting based on feature extraction and improved general regression neural network model’, *Energy* **166**, 653–663.
- [19] Luketina, J., Berglund, M., Greff, K. and Raiko, T. [2016], ‘Scalable gradient-based tuning of continuous regularization hyperparameters’.  
**URL:** <https://arxiv.org/abs/1511.06727>
- [20] Parkes, A. I., Savasta, T. D., Sobey, A. J. and Hudson, D. A. [2019], Efficient vessel power prediction in operational conditions using machine learning., in ‘Practical Design of Ships and Other Floating Structures(PRADS), September 2019, Yokohama, Japan’.
- [21] Parkes, A. I., Sobey, A. J. and Hudson, D. A. [2018], ‘Physics-based shaft power prediction for large merchant ships using neural networks’, *Ocean Engineering* **166**, 92–104.
- [22] Parkes, A. I., Sobey, A. J. and Hudson, D. A. [2021], ‘Towards error measures which influence a learners inductive bias to the ground truth’.  
**URL:** <https://arxiv.org/abs/2105.01567>

- [23] Pedersen, B. P. and Larsen, J. [2009], Prediction of full-scale propulsion power using artificial neural networks, *in* ‘Proceedings of the 8th international conference on computer and IT applications in the maritime industries (COMPIT’09), Budapest, Hungary May’, pp. 10–12.
- [24] Petersen, J. P., Jacobsen, D. J. and Winther, O. [2012], ‘Statistical modelling for ship propulsion efficiency’, *Journal of marine science and technology* **17**(1), 30–39.
- [25] Smith, C. and Jin, Y. [2014], ‘Evolutionary multi-objective generation of recurrent neural network ensembles for time series prediction’, *Neurocomputing* **143**, 302–311.
- [26] Tani, L., Rand, D., Veelken, C. and Kadastik, M. [2021], ‘Evolutionary algorithms for hyperparameter optimization in machine learning for application in high energy physics’, *The European Physical Journal C* **81**(2), 1–9.
- [27] Wang, B., Sun, Y., Xue, B. and Zhang, M. [2019], Evolving deep neural networks by multi-objective particle swarm optimization for image classification, *in* ‘Proceedings of the Genetic and Evolutionary Computation Conference’, GECCO ’19, p. 490–498.
- [28] Willard, J., Jia, X., Xu, S., Steinbach, M. and Kumar, V. [2020], ‘Integrating physics-based modeling with machine learning: a survey’.  
**URL:** <https://arxiv.org/pdf/2112.12979>
- [29] Yang, S., Tian, Y., He, C., Zhang, X., Tan, K. C. and Jin, Y. [2021], ‘A gradient-guided evolutionary approach to training deep neural networks’, *IEEE Transactions on Neural Networks and Learning Systems* pp. 1–15.
- [30] Zeiler, M. D. [2012], ‘Adadelta: an adaptive learning rate method’.  
**URL:** <https://arxiv.org/abs/1212.5701>